**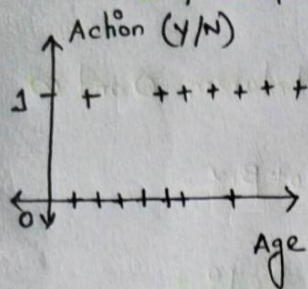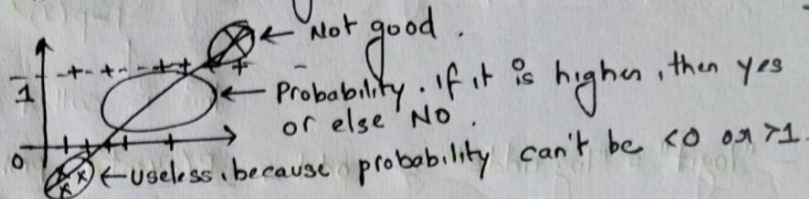Intution** → Suppose we send a email to a set of customers, whether they respond to the email or not. (Action - Yes/No).

Action (Y/N)

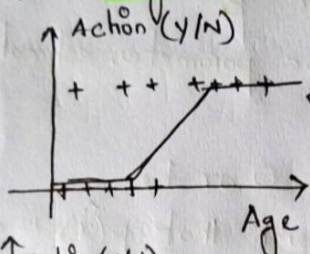→ We can observe, more the age they respond to the email.

→ It is not linear regression model because.

← Not good.

← Probability. If it is higher, then yes or else No.

← Useless, because probability can't be $<0$ or $>1$.

- We know probability $0-1$.
- So, the target variable is either yes/No. So that is classification.
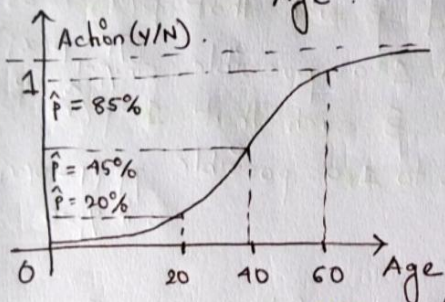
Action (Y/N)

best fit line.

○ Linear regression, $y = b_0 + b_1 * x$.

↓ Add Sigmoid function.

$$P = \frac{1}{1+e^{-y}}$$

↓

$$\ln\left(\frac{P}{1-P}\right) = b_0 + b_1 x \quad \text{(logistic equation)}$$

Action (Y/N)

$\hat{P} = 85\%$

$\hat{P} = 45\%$

$\hat{P} = 20\%$

0    20    40    60    Age

Mostly likely to respond - 60
Least likely to respond - 20.

if $\hat{P}$ is $> 50$ then yes for email.
$\hat{P}$ is $=< 50$ then No for email.

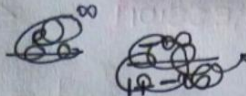[this number can change according to need].

Why logistic Regression?

→ Consider a dataset where response variable falls into one of two categories Yes or No. Rather than modelling the response Y directly, logistic regression models the probability that Y belong to particular category.

- For making a range of 0 to 1, we use logistic function, $P(x) = \dfrac{e^{b_0 + B_1 x}}{1 + e^{b_0 + B_1 x}}$

Method is known as maximum likelihood function.

- Range is 0 to 1. Values cannot be less than 0 or more than 1.
- Logistic function will always produce S curve regardless of value of x.

$$p(x) = \frac{e^{B_0 + B_1 x}}{1 + e^{B_0 + B_1 x}}, \quad \text{ececece etc.}$$

After bit of manipolation, $\dfrac{p(x)}{1 - p(x)} = e^{B_0 + B_1 x}$.

Quantity, $\dfrac{p(x)}{1 - p(x)}$ = odds, can take any value between $0$ and $\infty$.

If we apply log then $\log\left(\dfrac{p(x)}{1 - p(x)}\right) = B_0 + B_1 x$.

this is called logit function.

- "logit" = "log odds" $\Rightarrow$ odd $= \dfrac{P(event)}{1 - P(event)}$   probability event occur / probability event not occur.

- Parameter Estimation –
  - Goal of learning is to estimate parameter vector $\hat{B}$.
  - Logistic regression uses Maximum likelihood for parameter estimation.
  - How does Maximum likelihood works?
    
    → Consider N samples with labels either $0$ or $1$.
    
    → For sampled labelled "1": Estimated $\hat{B}$ such that $p(\hat{x})$ is as close to $1$ as possible. $\prod_{\sin y_i = 1} p(x_i)$
    
    → For sample labelled "0": Estimate $\hat{B}$ such that $1 - p(\hat{x})$ is as close to $1$ as possible. $\prod_{\sin y_i = 0}(1 - p(x_i))$

- We have to optimise likelihood function.

- If data is not linearly seperable, we cannot apply logistic regression. Linearly seperable is property of two sets of points. Two sets of points are linearly seperable if there exists atleast one line in the plane which seperate the both set of points (blue and red).

Main steps in logistic Regression – Training set is given.

Step 1 – How to calculate logistic function. (to learn paramoter of training set)
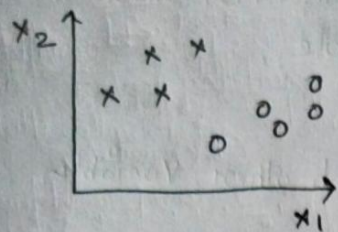
Step 2 – How to learn the coefficient for a logistic regression model using stochastic gradient.

Step 3 – How to make prediction using the model.

Logistic fn – logit fn / sigmoid fn to translate value of X to graph & learn the coefficient ($B_0 / B_1 / B_2$ etc). Build the model and predict.
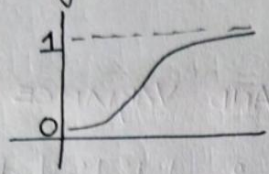
Let's consider



Step 1a - Linear seperable.

Step 1b - We need to transform data point using logit $f^n$ or Sigmoid $f^n$ which is given by

$$\underset{\text{(logit } f^n)}{h_\theta(x)} = \frac{1}{1+e^{-x}} \quad —① $$

Function transform each input value in range 0 to 1.

Step 2a - After transformation, -ve number resulted in value close to zero. The larger +ve number resulted in value close to one.
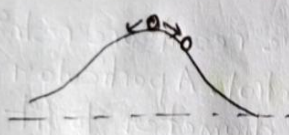


$$P_{(class = 0)} = \frac{1}{1+e^{-(b_0 + b_1 x_1 + b_2 x_2)}}$$

class = 0 or 1 . modified logit $f^n$.

if Probability $\geq 0.5$ then 1 or else 0.

Step 2b - We need to find parameter of eqn $(b_0, b_1, b_2)$ with the help of gradient descent (to find best value of parameter).

Gradient descent → ① Suppose you are in top of mountain. & choose the best route to reach ground level.
② Take a step in any direction.
③ Calculate distance at all angle and choose the smallest distance angle and take step.
④ Repeat step ③ till you reach the ground.

so, for first step choose 0 for parameter.

$$\text{prediction} = \frac{1}{1+e^{-(0 + 0(2.7) + 0(2.3))}}$$

$= 0.5$ ell

Ultimate we have to get best value (so that it will be equal to predicted) In short accuracy & minimum error.

$-\alpha$ is a learning rate, at which model can learn.

Q) What is maximum likelihood function?

- Goal of the Maximum likelihood estimation is to make inferences about the population that is most likely to generate a sample.

Q) Assumptions in logistic regression?

① LR requires observation to be independent of each other. Variables should not be highly correlated.

② Little or no multicollinearity among independent variables

③ LR assum linearity of independent variable and log odds.

④ LR typically require large sample size.

Topic o  ⚬ o Cost

# INTERPRETATION OF STANDARD DEVIATION AND VARIANCE

- Standard deviation = $\sqrt{Variance}$

→ Basically a small SD means that the values in statistical dataset are close to the mean of the dataset, on average and a large standard deviation means that the value in the dataset are farther away from the mean, on average.

→ In short, it measure how concentrated the data around mean, more concentrated → smaller SD.

→ A small SD can be the goal in certain situations where result are restricted for example in product manufacturing and quality control. A particular type of car part that has to be 2 centimeters in diameter to fit properly had better not to have a very big standard deviation during manufacturing process. A big standard deviation would mean it will end up in trash can because they don't fight right.

→ High SD, reflects a large amount of variation in the group. For example, if we look at the salaries for everyone in a company, including student intern to the CEO, standard deviation may be very large. On other hand, if we observe only student interns salary, standard deviation may be low/smaller.

→ Outlier does affect the SD, because formula includes the mean. SD cannot be negative and lowest possible value is 0. and 0 is possible only when every single entity have same number (no deviation)

→ SD have the same unit as the original data.